

I hereby certify that this paper is being deposited with the United States Postal Service as Express Mail in an envelope addressed to: U.S. Patent and Trademark Office, Washington, D.C. 22202 on this date.

November 21, 2001

Date

Express Mail Label No.:

EL84622235US

METHODS AND APPARATUS FOR GENERATING MULTIPLE NETWORK STREAMS FROM A LARGE SCALE MEMORY BUFFER

09990093-12101

This invention relates to methods and apparatus for attaching multiple network processors, or output stream processors, to a large-scale unitary memory array, and more particularly, to apparatus for streaming electronic information from a large-scale unitary memory array to multiple attached networks, using hardware-based network processors or stream processors that generate transport protocols for the streams. A hardware-based arbitrator controls access of the stream processors to the memory array.

BACKGROUND OF THE INVENTION

The ability to share a large memory array between multiple network processors allows multiple output streams to be generated from a single copy of material, and to be transmitted simultaneously through many output ports at different intervals, without having to replicate or make copies of the content. The larger the memory buffer is, the more data or unique programs or source material can be stored. By utilizing a unitary memory array, it is possible to generate many outputs from a single copy of the program material. The present invention is especially well suited to

1 audio and video streaming applications where the program material size, or payload, is
2 large, and a great number of simultaneous playback streams need to be generated.

3 When building audio, video, or Internet streaming servers (collectively
4 "servers"), there are great demands placed on the architecture. There are several
5 design parameters that must be considered. One is the number of simultaneous output
6 streams or playback sessions that are required. Another is the total size in bytes of the
7 information to be streamed. A third parameter is the data rate, or bit-rate of each
8 output stream. By knowing these parameters, one can specify a target implementation
9 architecture.

10 There are many types of existing high-speed network interfaces to
11 implement the output connection. Typically, these exist as network controller chips.
12 These network controllers are well suited for central processing unit (CPU) based
13 computers, such as a PC or Workstation. PC's and workstations usually contain a PCI
14 bus, or similar expansion interface for adding network and I/O controller cards. Such
15 expansion buses are designed for a single processor as might be found in a desktop
16 computer system. Moreover, computers contain only a single network controller. In
17 specific cases, such as a file server, there may be multiple network interfaces, but
18 usually no more than two.

19 An object of a streaming server is to stream as much content as possible
20 to as many users as possible, from as small a space as possible. However, there are
21 limitations on the number of network interfaces that can be added to a computer
22 towards this goal.

Typical computers contain a single CPU or central processing unit. Among other things, the CPU is responsible for running software to generate all of the network traffic, usually called a transport protocol stack. However, the speed of the CPU then becomes a limitation or performance bottleneck. To overcome this limitation, multiple CPU's are usually added. However, implementation of this approach for the present application would require multiple network interface cards, and multiple CPU's all competing for the interconnect or data bus structure in the server. Some solutions to this bottleneck have been devised, with varying levels of success.

Fig. 9 shows one implementation of the prior art. In this configuration, a CPU 902 controls and implements all data transfers. A data block is retrieved from a storage device 901 by the CPU 902 over signal lines 911, and the data block is subsequently written to a memory 903 over signal lines 912. After a complete block is stored in the memory 903, the CPU 902 can generate appropriate networking protocol packets and store them in the memory 903, under software control. Once a protocol packet has been stored in the memory 903, the CPU 902 can move the packet to the output interface 904, over signal lines 914. The data block is sent to a client device through line 915.

The final parameter relates to storage. All of the streaming payload data must originate in the server. In current architectures, the data is usually stored in hard disk drives. All hard disk drives have two main limiting factors, the sustained transfer rate of the data, and the seek time or random access time. The transfer rate dictates

1 how fast data can be read from the drive in a continuous manner. Seek time dictates
2 the latency for moving from one part of the drive to another to access a specific data
3 block. Any time spent seeking to a specific location on a disk takes away from the
4 time that could be used for transferring data. In this way, the efficiency of a disk drive
5 can be greatly reduced by the seek time. In a streaming server application, there can
6 be many output streams, each representing different programs, or different locations in
7 the same program. This creates a bottleneck in the storage subsystem, as the disk drive
8 will spend a significant amount of time seeking for the appropriate data block. To
9 overcome this, more hard drives are added.

10 The greater the number of output streams, the more hard drives the
11 system will require. Hard drives can be combined into arrays or groups of drives,
12 such as RAID (Redundant Array of Inexpensive Disks) and JBOD (Just a Bunch Of
13 Disks) configurations, but ultimately, all rotational media has transfer rate and seek
14 time limitations.

15 Data blocks read from the storage device 901 can be of different size
16 than the final protocol block size. This buffer size translation can be accomplished
17 under program control by the CPU 902. This system is economical since there is no
18 substantial hardware, and all functions are accomplished by the software on the CPU
19 902. The primary limitation to this approach is performance, which is constrained by
20 the CPU processing time.

21 Fig. 10 is an improved version of the prior art just described that
22 implements Direct Memory Access or DMA. In this version, a shared bus 1011 is

added to allow a storage device 1001 and an output interface 1004 to directly access the memory 1003. The CPU 1002 begins by setting up a transfer with storage device 1001, by using signal lines 1013, bus 1011, and signal lines 1012. The storage device 1001 then begins initiating a transfer to memory 1003, over signal lines 1012, bus 1011, and signal lines 1014, to memory 1003. The transfer can occur without the CPU 1002 being in the data path, which increases performance.

Once a block of data is in memory 1003, the CPU 1002 generates the appropriate networking protocol packets and stores them in memory 1003, under software control. Once a protocol packet has been stored in the memory 1003, the CPU 1002 sets up the output transfer to the output interface 1004. The output interface 1004 then initiates a transfer from memory 1003, over signal lines 1014, bus 1011, signal lines 1015, and through the output interface 1004 to the output 1016. In this system, the CPU is not responsible for actually moving the data, which increases performance when compared to the system in Fig. 9. However, the protocol packets are still generated by the CPU 1002, the memory 1003 has a relatively small size, and the bus 1011 must be shared with all devices. Even with the fastest bus and the fastest CPU, this architecture is limited in capacity when compared to the inventive system.

Accordingly, one object of the present invention is to increase the number of output streams possible from a single memory. Another object of the invention is to increase the size of the data bus and address bus so that a higher data rate can be achieved from a much larger memory array. Another object of the

1 invention is to remove the generic CPU from the memory bus arbitration process and
2 the protocol stack generation process.

SUMMARY OF THE INVENTION

3 In keeping with one aspect of this invention, a large scale memory stores
4 a number of video, audio, audiovisual and other content. The memory is random
5 access, eliminating the access times required by hard disk drives. The content can be
6 read out of the memory to multiple customer sites over several networks, many of
7 which use different network protocols. Content is stored in the memory and read out
8 of the memory to the various customer sites under the control of a hardware based
9 arbitrator.

10 The content is bundled into data packets, each of which is encoded with
11 a protocol to form a transport protocol stack. The transport protocol stack is generated
12 in hardware-based architecture, thus greatly improving the throughput, and increasing
13 the number of streams that can be generated. A wide data bus and wide address bus
14 can be utilized because the protocol stack is generated in hardware, so that higher
15 throughput can be achieved from a large scale memory. A plurality of protocol stack
16 generators have access to the same block of memory, allowing many output streams to
17 be generated from a single copy of content in the large scale memory.

BRIEF DESCRIPTION OF THE DRAWINGS

18 The above mentioned and other features of this invention and the
19 manner of obtaining them will become more apparent, and the invention itself will be
20

best understood by reference to the following description of an embodiment of the invention taken in conjunction with the accompanying drawings, in which:

FIG. 1 is a block diagram of an embodiment of a communication system made in accordance with the present invention;

FIG. 2 is a block diagram of a portion of the system of Fig. 1;

FIG. 3 is a flow chart showing the operation of the arbitrator of Fig. 2;

FIG. 4 is a block diagram of a stream server module in the system of Fig. 1;

FIG. 5 is a block diagram of the stream server processor in the stream server module of Fig. 4;

FIG. 6 is a block diagram showing the input and output paths of a stream controller in the stream server processor of Fig. 5;

FIG. 7 is a state diagram showing the operation of the stream controller of Fig. 6;

FIG. 8 is a flowchart of the operation of the protocol encoder logic of Fig. 5, referenced in state S704 in Fig. 7;

FIG. 9 is a block diagram of a conventional communication system; and

FIG. 10 is a block diagram of another conventional communication system.

DETAILED DESCRIPTION

As seen in Fig. 1, a server system 100 is primarily built from a memory array 101, an interconnect device 102, and stream server modules 103a through 103n (103). The server system 100 is part of a communication system that also includes transport networks 122a through 122n (122), and client devices 124a through 124n (124). In a typical system, each client device would operate through a single transport network, but each transport network could communicate with the server system 100 through any of the stream server modules.

Each transport network can operate using a different protocol, such as IP (Internet Protocol), ATM (Asynchronous Transfer Mode), Ethernet, or other suitable Layer-2 or Layer-3 protocol. In addition, a specific transport network can operate with multiple upper-level protocols such as Quick Time, Real Networks, RTP (Real Time Protocol), RTSP (Real Time Streaming Protocol), UDP (User Datagram Protocol), TCP (Transport Control Protocol), etc. A typical example would be an Ethernet transport network with IP protocol packets that contain UDP packets, which in turn contain RTP payload packets.

The communication process starts with a stream request being sent from a client device 124 over an associated transport network 122. The command for the request arrives over a signal line 114a-114n (114) to a stream server module 103, where the protocol information is decoded. If the request comes in from stream server module 103a, for example, it travels over a bus 117 to a master CPU 107. The master

1 CPU can be implemented using any number of commercially available CPU's, one
2 such CPU being a PowerPC 750 made by Motorola. For local configuration and
3 status updates, the CPU 107 is also connected to a local control interface 106 over
4 signal line 120, which communicates with the system operator over a line 121.
5 Typically this could be a terminal or local computer using a serial connection or
6 network connection.

7 Control functions, or non-streaming payloads, are handled by the master
8 CPU 107. Program instructions in the master CPU 107 determine the location of the
9 desired content or program material in memory array 101. The memory array 101 is a
10 large scale memory buffer that can store video, audio and other information. In this
11 manner, the server system 100 can provide a variety of content to several customer
12 devices simultaneously. Customer sessions can include movies, music, sports events,
13 written information, etc., each of which can represent a program. However, each
14 customer device can receive the same content or different content. Each customer
15 receives a unique asynchronous stream of data that might or might not coincide in
16 time with unique asynchronous streams sent to other customer devices.

17 If the requested content is not already resident in the memory array 101,
18 a request to load the program is issued over signal line 118, through a backplane
19 interface 105 and over a signal line 119. An external processor or CPU (not shown)
20 responds to the request by loading the requested program content over a backplane
21 line 116, under the control of backplane interface 104. Backplane interface 104 is
22 connected to the memory array 101 through the interconnect 102. This allows the

1 memory array 101 to be shared by the stream server modules 103, as well as the
2 backplane interface 104. The program content is written from the backplane interface
3 104, sent over signal line 115, through interconnect 102, over signal line 112, and
4 finally to the memory array 101.

5 Backplanes typically operate more efficiently when moving data in
6 chunks, or blocks. As such, backplane interface 104, interconnect 102, and memory
7 array 101 can each contain small buffers to allow larger 'bursts' of data to be
8 transferred. Another way to achieve higher speeds is to use a wider bus path, such as
9 128 bits, 256 bits, or larger. A wider bus interface allows more bytes of data to be
10 transferred on each memory access cycle.

11 When the first block of program material has been loaded into memory
12 array 101, the streaming output can begin. Streaming output can also be delayed until
13 the entire program has been loaded into memory array 101, or at any point in between.
14 Data playback is controlled by a selected one or more stream server modules 103. If
15 the stream server module 103a is selected, for example, the stream server module
16 103a sends read requests over signal line 113a, through the interconnect 102, over a
17 signal line 111 to the memory array 101. A block of data is read from the memory
18 array 101, sent over signal line 112, through the interconnect 102, and over signal line
19 113a to the stream server module 103a. Once the block of data has arrived at the
20 stream server module 103a, the transport protocol stack is generated for this block and
21 the result is sent to transport network 122a over signal line 114a. Transport network
22 122a then carries the steaming output data block to the client device 124a over signal

0900003-12107

1 line 123a. This process is repeated for each data block contained in the program
2 source material.

3 If the requested program content already resides in the memory array
4 101, the CPU 107 informs the stream server controller 103a of the actual location in
5 the memory array. With this information, the stream server module can begin
6 requesting the program stream from memory array 101 immediately.

7 The system is broken into two separate paths; the first is for large
8 content or payload; the second is for control and other non-payload types of packets.
9 Non-payload packets could be VCR type controls or "Trick Mode" packets, such as
10 Pause, Fast-Forward, Rewind, etc., as well as program listings, or content availability
11 information. Since these signals are generally not very CPU demanding, they can be
12 easily handled by a CPU running software. The actual payload packets are very CPU
13 intensive, yet little processing needs to be done to the payload. In this case, the stream
14 server module 103a, or other stream server module, can handle the transfer and
15 movement of payload data, without requiring participation by the master CPU 107.
16 This separation of the paths allows a much higher system density, and provides a
17 much greater stream capacity when compared to CPU based designs.

18 Memory array 101 is preferably large enough to hold many source
19 programs, and could be many Gigabytes or Terabytes in size. A minimum of 65
20 gigabytes is recommended. Such source programs can be in the form of video, audio,
21 or data, including but not limited to any combination thereof. Such memory arrays
22 may be built from conventional memory including but not limited to dynamic random

access memory (DRAM), synchronous DRAM (SDRAM), Rambus DRAM (RDRAM), dual data rate DRAM (DDRDRAM), static RAM (SRAM), magnetic RAM (MRAM), flash memory, or any memory that is solid state in nature. Dual inline memory modules (DIMMs) can be used. In order to access such a large memory array, a wide address bus is used. A conventional 32-bit address bus is only sufficient to address 4 Gigabytes of RAM, and is not preferred for this application. An address bus greater than 36 bits wide is preferred, and a 48 bit address bus would be more suitable for this application because it can directly access 256 Petabytes of memory.

The interconnect 102 is shown in greater detail in Fig. 2. The interconnect 102 controls the transfer of data between the memory array 101 and the stream server modules 103. The interconnect 102 also establishes priority among the stream server modules 103, determining the order in which the stream server modules receive data from the memory 101.

The interconnect 102 includes an arbitrator 202. The arbitrator 202 is preferably hardware based, and can include a field programmable gate array or another suitable device. The arbitrator 202 is a hardware based state machine. In prioritizing requests, the arbitrator 202 could be programmed to give the backplane interface 104 top priority, if desired. The several stream server modules 103 can be given priority in any suitable manner, such as serial priority, priority based on content such as audio, video, etc. or any other desirable manner.

Each stream server module is connected to an address bus 111 through signal lines 208a...208n (208). Data is sent from the memory array 101 to the stream server modules 103 over a data bus 112 and signal lines 210a-210n (210).

The stream server modules request data through the arbitrator 202. For example, the stream server module 103a sends requests for data to the arbitrator 202 over the signal line 204a. When the arbitrator decides that the stream server module 103a should receive data, an authorization is sent over a line 210a. In this manner, the arbitrator sets priority with respect to the stream server modules.

The backplane interface 104 is used to load new information into the memory array 101, among other things. The new information is provided through the signal line 116 through the backplane interface 104. When the backplane interface 104 receives data, it requests access to the address bus 111 and the data bus 112 from the arbitrator 202 through a signal line 212. The arbitrator 202 authorizes data transfer over a signal line 214. Upon authorization, the backplane interface 104 provides address data to the memory array 101 over the bus 111, using the signal line 216. Data is transferred over the bus 112 through a signal line 218.

The operation of the arbitrator 202 is shown in greater detail in the flow chart of Fig. 3. At step S302, the arbitrator 202 determines what, if any, request signals have been received from the stream server modules 103 and the backplane interface 104. If there are no requests, the arbitrator waits for requests to be received by looping around steps S302 and S304. If one or more requests have been received at step S304, then the arbitrator stores all requesting devices, resets the selected

“winning” device, and sets a pointer to the first requesting device at step S306. Since there can be multiple devices all requesting simultaneously, the highest priority device must be selected as the “winner”. S308 checks to see if the currently selected requestor has the highest priority. If it does, then the new “winner” is selected in S310. The process continues for all requesting devices through S312, which checks for the last device. If this is not the last one, S314 will increment the pointer to select the next device. The process repeats until all requests have been evaluated for their priority, the highest priority device being designated as the “Winner”. Once complete, control is granted to the “Winner” in S316. The winner can now issue commands through the arbitrator. Service remains granted to this “winner” as long as the currently selected device demands service. S318 monitors the request signal from the arbitrated winner, and holds the control grant as long as the request is present by looping through S318 and S316. Once the device has stopped requesting service, as determined in S318, S320 releases control from the winning requestor, and passes control back to S302, which starts the arbitration again.

Fig. 4 is a block diagram of an implementation of the stream server modules 103 shown in Fig. 1. A stream server processor (SSP) 401 serves as the automatic payload requester, as well as the protocol encoder and decoder. The SSP 401 requests and receives data payload over signal line 113. It then encodes and forms network packets, such as TCP/IP or UDP/IP or the like. The encoded packets are sent out over signal lines 411a-411n (411), to one or more media access controllers (MAC) 402a-402n (402). The media access controllers 402 handle the

1 serialization and de-serialization of data and formatting as required by the specific
2 physical network used. In the case of Ethernet, the Media Access Controllers 402 also
3 handle the detection of collisions and the auto-recovery of link-level network errors.

4 The media access controllers 402 are connected utilizing signal lines
5 412a-412n (412), to media interface modules 403a-403n (403), which are responsible
6 for the physical media of the network connection. This could be a twisted-pair
7 transceiver for Ethernet, Fiber-Optic interface for Ethernet, SONET or many other
8 suitable physical interfaces, which exist now or will be created in the future, such
9 interfaces being appropriate for the physical low-level interface of the desired
10 network, and sent out over the signal lines 114a-114n (114).

11 When control packets are required, such as VCR like controls, it is more
12 efficient to handle the processing in a CPU instead of hardware. Depending on the
13 installation requirements, the protocols can change frequently, and could vary from
14 system to system. In certain cases, the control protocol stack could be customized
15 depending on the installation environment. For these reasons, the control functions
16 can be implemented in CPU 404. By contrast, however, the actual payload protocol
17 stack is relatively stable, so it can be processed in hardware. Additionally, a hardware
18 implementation allows for a wide data path, much wider than the data path for a
19 standard CPU, which typically has a 32 bit or 64 bit data path. By using a wider bus,
20 such as a 256 bit data bus, much more data can be moved from memory on each
21 access cycle.

1 onto line 517a, into the control data buffer 504a, over line 522a, and into the protocol
2 stream encoder/decoder 505a, where they are encoded, and are then sent out over the
3 line 411a. Buffering the control data in the control data buffer 504a allows the control
4 data bus interface 510 and associated busses to operate at a separate transfer speed
5 from that of the protocol stream encoder/decoder 505a.

6 Fig. 6 is a detailed block diagram of a portion of one of the stream
7 controllers 501. Outgoing payload data is provided on a line 516 and is clocked
8 through the payload data buffer 503. The data is sent over line 611, to a protocol
9 select logic array 601, which sends the data to an appropriate protocol encoder logic
10 array 602a-602n (602), through lines 612. After the data blocks are encoded with the
11 correct protocol, they are sent over lines 613a-613n (613) to the payload/control
12 transmit combiner logic 403, and to the network interface over line 411.

13 Outgoing control data for a particular transmission is sent over line 517
14 to a control data buffer 504. The output of the control data buffer 504 is sent over line
15 614 to protocol select logic circuitry 605, which identifies the required protocol, and
16 sends the data over one of the lines 615 to an appropriate protocol encoder logic
17 circuit 604a-604n (604). The encoded data is sent to the payload/control transmit
18 combiner logic 603 over lines 616a-616n (616), and the control information is
19 transmitted over line 411.

20 Incoming data is decoded in a similar manner. Data entering the system
21 through line 411 is stripped in payload/control receive separator logic 606. The
22 payload data is sent over line 617 to protocol select logic 607. The protocol select

1 logic identifies the protocol used for the particular payload, and sends the data over an
2 appropriate line 618 to the correct protocol decoder logic circuitry 608a-608n (608).
3 The output of the decoded data is clocked through payload data buffer 503 over lines
4 619, and sent out over line 516.

5 Control data is decoded in a similar manner. Entering the system
6 through line 411, the data is stripped in payload/control receive separator logic 606.
7 The stripped data is sent over line 620 through protocol select logic 609, which
8 identifies the protocol used for the control data, and sends the data over an appropriate
9 line 621 to the desired protocol decoder logic 610a-610n (610). The decoded data is
10 sent over line 622 and clocked through a selected control data buffer 504, leaving the
11 stream controller through line 517.

12 Fig. 7 is a state diagram for the stream controllers 501 shown in Fig. 5.
13 The stream controller operates in an idle state S701 until a block of content data is
14 required. At that time, the appropriate addresses are generated at S702, and the
15 content data is read from the memory 101 as a burst in state S703. The memory 101
16 burst is read until the buffer is full. Then, if needed, time stamps are adjusted and a
17 protocol wrapper is added at state S704 until all the data in the buffer is encoded with
18 the appropriate protocol. The operational state details of S704 are further described in
19 Fig. 8. The data in the buffer is written to an output interface in state S705, and the
20 stream controller returns to the idle state at S701.

21 Fig. 8 details the operation of the protocol encoder logic 505 as shown
22 in Fig. 5, which is referenced in state S704 of Fig. 7. At step S801 the buffer pointer

